# CS 595: Assignment #8

Due on Friday, November 14, 2014

*Dr Nelson 4:20pm*

**Victor Nwala**

# Contents

# Problem 1

1. What 5 movies have the highest average ratings? Show the movies and their ratings sorted by their average ratings.

The code sorts the result from least to highest rated movie. Their rankings are:

1) 1201—Marlene Dietrich: Shadow and Light (1996) 2) 1293—Star Kid (1997) 3) 1189—Prefontaine (1997) 4) 1653—Entertaining Angels: The Dorothy Day Story (1996) 5) 1536—Aiqing wansui (1994) All with average ratings of 5.0

Listing 1: Python script solving problem 1

```python
from collections import Counter
from collections import defaultdict
from operator import itemgetter, attrgetter
from collections import OrderedDict


path='/home/vnwala/ml-100k'
def getMovieData():
    movie = {}
    for line in open(path + '/u.data'):
        (user, movieid, rating, ts) = line.split('\t')

        if( movieid in movie ):
            movie[movieid].append(float(rating))
        else:
            movie[movieid] = []
            movie[movieid].append(float(rating))

    return movie


movie = getMovieData()

sums = {}

for movieid in movie:
    a = sum(movie[movieid])/len(movie[movieid])
    if( movieid in movie ):
        sums[movieid] = a
    else:
        sums[movieid] = a

sums = sorted(sums.items(), key=lambda x: x[1])
print sums
```
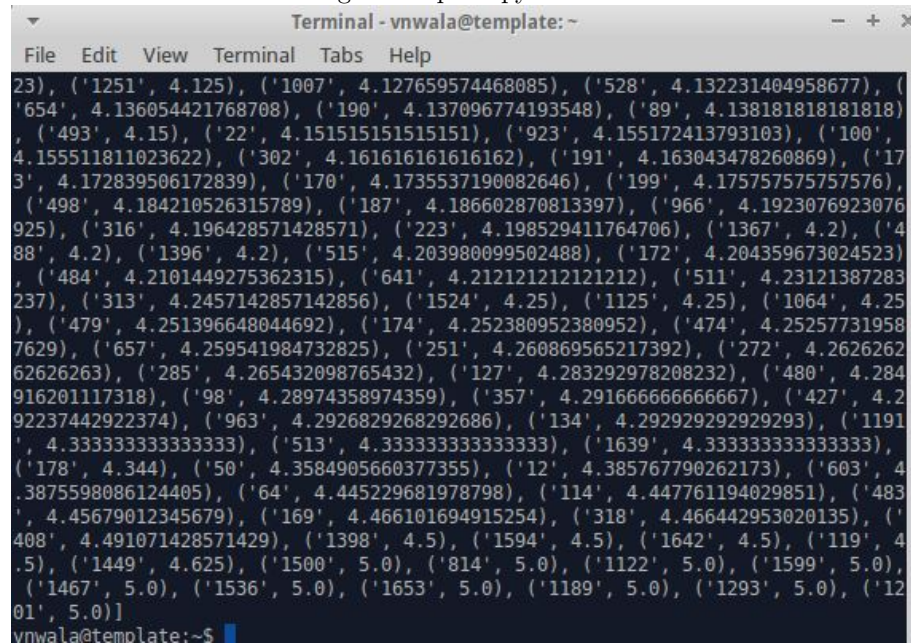
Figure 1: ques1.py at work

# Problem 2

2. What 5 movies received the most ratings? Show the movies and the number of ratings sorted by number of ratings.

The code sorts the result from least to most rated movie. Their rankings are:

1) 50—Star Wars (1977) with 583 ratings 2) 258—Contact (1997) with 509 ratings 3) 100—Fargo (1996) with 508 ratings 4) 181—Return of the Jedi (1983) with 507 ratings 5) 294—Liar Liar (1997) with 485 ratings

Listing 2: Python script solving problem 2

```python
from collections import Counter
from collections import defaultdict
from operator import itemgetter, attrgetter
from collections import OrderedDict

path='/home/vnwala/ml-100k'
def getMovieData():
    movie = {}
    for line in open(path + '/u.data'):
        (user, movieid, rating, ts) = line.split('\t')

        if( movieid in movie ):
            movie[movieid].append(float(rating))
        else:
            movie[movieid] = []
            movie[movieid].append(float(rating))

    return movie


movie = getMovieData()

sums2 = {}

for movieid in movie:
    a = len(movie[movieid])
    if( movieid in movie ):
        sums2[movieid] = a
    else:
        sums2[movieid] = a
sums2 = sorted(sums2.items(), key=lambda x: x[1])
print sums2
```
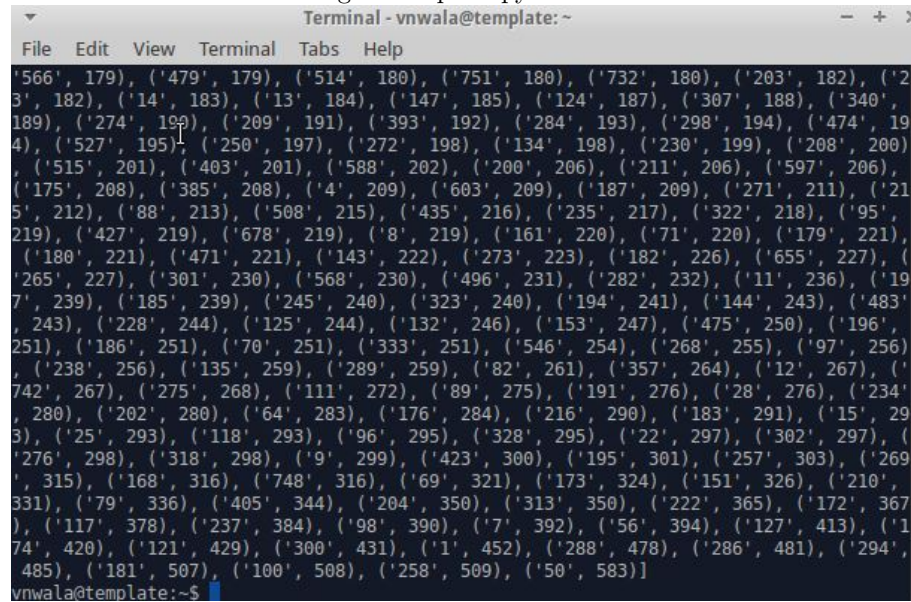
      

Figure 2: ques2.py at work

# Problem 3

3. What 5 movies were rated the highest on average by women? Show the movies and their ratings sorted by ratings.
1) 286—English Patient, The (1996) with 152 female raters 2) 50—Star Wars (1977) with 151 female raters 3) 288—Scream (1996) with 143 female raters 4) 294—Liar Liar (1997) with 141 female raters 5) 258—Contact (1997) with 137 female raters

Listing 3: Python script solving problem 3

```python
path='/home/vnwala/ml-100k'
user1 = {}
movie = {}
female_raters ={}

for line in open(path + '/u.user'):
    (userid , age, gender, occupation, zipcode) = line.split('|')
    if( gender ==  'F' ):
        user1[userid] = (age, gender, occupation, zipcode)



        for line in open(path + '/u.data'):
            (user, movieid, rating, ts) = line.split('\t')
            if (user == userid):
                if( movieid in movie ):
                    movie[movieid].append(userid)
                else:
                    movie[movieid] = []
                    movie[movieid].append(userid)
for movieid in movie:
    female_raters[movieid] = len(movie[movieid])

female_raters = sorted(female_raters.items(), key=lambda x: x[1])

print female_raters
```

Figure 3: ques3.py at work

# Problem 4

4. What 5 movies were rated the highest on average by men? Show the movies and their ratings sorted by ratings.

1) 50—Star Wars (1977) with 432 male raters 2) 181—Return of the Jedi (1983) with 383 male raters 3) 100—Fargo (1996) with 383 male raters 4) 258—Contact (1997) with 372 male raters 5) 294—Liar Liar (1997) with 344 male raters

Listing 4: Python script solving problem 4

```
path='/home/vnwala/ml-100k'
user1 = {}
movie = {}
male_raters ={}

for line in open(path + '/u.user'):
    (userid , age, gender, occupation, zipcode) = line.split('|')
    if( gender ==  'M' ):
        user1[userid] = (age, gender, occupation, zipcode)



        for line in open(path + '/u.data'):
            (user, movieid, rating, ts) = line.split('\t')
            if (user == userid):
                if( movieid in movie ):
                    movie[movieid].append(userid)
                else:
                    movie[movieid] = []
                    movie[movieid].append(userid)
for movieid in movie:
    male_raters[movieid] = len(movie[movieid])

male_raters = sorted(male_raters.items(), key=lambda x: x[1])

print male_raters
```

     

Figure 4: ques4.py at work

# Problem 5

5. What movie received ratings most like Top Gun? Which movie received ratings that were least like Top Gun (negative correlation)?

Some movies rated the most like like Top Gun are:

1) Mr. Smith Goes to Washington (1939) 2) (Cold Fever) (1994) 3) Young Guns II (1990) 4) Young Poisoner's Handbook, The (1995) 6) Zeus and Roxanne (1997) 7) Young Poisoner's Handbook, The (1995) 8) Out to Sea (1997) 9) Old Yeller (1957) 10) Hungarian Fairy Tale, A (1987)

Listing 5: Python script solving problem 5

```python
from collections import Counter
from collections import defaultdict
from operator import itemgetter, attrgetter
from collections import OrderedDict

path='/home/vnwala/ml-100k'




path='/home/vnwala/ml-100k'


def sim_distance(prefs, p1, p2):
    '''
    Returns a distance-based similarity score for person1 and person2.
    '''
    # Get the list of shared_items
    si = {}
    for item in prefs[p1]:
        if item in prefs[p2]:
            si[item] = 1
    # If they have no ratings in common, return 0
    if len(si) == 0:
        return 0
    # Add up the squares of all the differences
    sum_of_squares = sum([pow(prefs[p1][item] - prefs[p2][item], 2) for item in
    prefs[p1] if item in prefs[p2]])
    return 1 / (1 + sum_of_squares)


def sim_pearson(prefs, p1, p2):
    '''
    Returns the Pearson correlation coefficient for p1 and p2.
    '''
    # Get the list of mutually rated items
    si = {}
    for item in prefs[p1]:
        if item in prefs[p2]:
            si[item] = 1
    # If they are no ratings in common, return 0
    if len(si) == 0:
        return 0
```

```
         # Sum calculations
45       n = len(si)
         # Sums of all the preferences
         sum1 = sum([prefs[p1][it] for it in si])
         sum2 = sum([prefs[p2][it] for it in si])
         # Sums of the squares
50       sum1Sq = sum([pow(prefs[p1][it], 2) for it in si])
         sum2Sq = sum([pow(prefs[p2][it], 2) for it in si])
         # Sum of the products
         pSum = sum([prefs[p1][it] * prefs[p2][it] for it in si])
         # Calculate r (Pearson score)
55       num = pSum - sum1 * sum2 / n
         den = sqrt((sum1Sq - pow(sum1, 2) / n) * (sum2Sq - pow(sum2, 2) / n))
         if den == 0:
             return 0
         r = num / den
60       return r


    def topMatches(
    prefs,
65  person,
    n=5,
    similarity=sim_pearson,
    ):
         '''
70       Returns the best matches for person from the prefs dictionary.
         Number of results and similarity function are optional params.
         '''
         scores = [(similarity(prefs, person, other), other) for other in prefs
         if other != person]
75       scores.sort()
         scores.reverse()
         return scores[0:n]




80




85  def transformPrefs(prefs):
         '''
         Transform the recommendations into a mapping where persons are described
         with interest scores for a given title e.g. {title: person} instead of
         {person: title}.
90       '''
         result = {}
         for person in prefs:
             for item in prefs[person]:
                 result.setdefault(item, {})
95       # Flip item and person
             result[item][person] = prefs[person][item]
```

```
            return result




    def calculateSimilarItems(prefs, n=5):
        '''
        Create a dictionary of items showing which other items they are
        most similar to.
        '''
        result = {}
        # Invert the preference matrix to be item-centric
        itemPrefs = transformPrefs(prefs)
        c = 0
        for item in itemPrefs:
        # Status updates for large datasets

            scores = topMatches(itemPrefs, 'Top Gun (1986)', n=n, similarity=
                sim_distance)
            result[item] = scores
        return result




    # Get movie titles
    movies = {}
    for line in open(path + '/u.item'):
        (id, title) = line.split('|')[0:2]
        movies[id] = title
    # Load data
    prefs = {}
    for line in open(path + '/u.data'):
        (user, movieid, rating, ts) = line.split('\t')
        prefs.setdefault(user, {})
        prefs[user][movies[movieid]] = float(rating)


    output = calculateSimilarItems(prefs)

    print output
```

Figure 5: ques5.py at work

# Problem 6

6. Which 5 raters rated the most films? Show the raters' IDs and the number of films each rated.

The solution is computed in this format, rater id:number of movies rated

1) 405:737 2) 655:685 3) 13:636 4) 450:540 5) 276:518

Listing 6: Python script solving problem 6

```python
from math import sqrt
import os, sys


def loadMovieLens(path='/home/vnwala/ml-100k/'):
    movies = {}
    for line in open(path + 'u.item'):
        (id, title) = line.split('|')[0:2]
        movies[id] = title
    prefs = {}
    for line in open(path + 'u.data'):
        (user, movieid, rating, ts) = line.split('\t')
        prefs.setdefault(user, {})
        prefs[user][movieid] = float(rating)
    return prefs, movies



def aggregateMovieAndUserData(path='/home/vnwala/ml-100k/'):

    try:

        movies = {}
        aggregateMovieData = []
        for line in open(path + 'u.item'):
            (id, title) = line.split('|')[0:2]
            movies[id] = (title, [], -1)

        users = {}

        for line in open(path + 'u.data'):
            (user, movieid, rating, ts) = line.split('\t')

            user = user.strip()
            movieid = movieid.strip()
            rating = rating.strip()
            ts = ts.strip()

            users.setdefault(user, {})
            users[user][movieid] = float(rating)


            movies[movieid][1].append(float(rating))

        for movieId, tupleData in movies.items():
            averageRating = sum(tupleData[1]) / float(len(tupleData[1]))
```

```
                    movietuples = (movieId, tupleData[0], averageRating, len(tupleData[1])
                        )
                    aggregateMovieData.append(movietuples)


50          aggregateUserData = []
            for line in open(path + 'u.user'):
                (userId, age, gender, occupation, zipCode) = line.split('|')
                userTuples = (userId, gender, age, users[userId])
                aggregateUserData.append(userTuples)
55

        except:
            exc_type, exc_obj, exc_tb = sys.exc_info()
            fname = os.path.split(exc_tb.tb_frame.f_code.co_filename)[1]
60          print(fname, exc_tb.tb_lineno, sys.exc_info() )
            return




65      return aggregateMovieData, aggregateUserData, movies


    def getHighestRatedMovies(movies, count):

70      if( count > 0 and count < len(movies) ):
            movies = sorted(movies, key=lambda tup: tup[2], reverse=True)
            i = 1
            for movie in movies:
                print movie
75
                if( i == count ):
                    break
                i = i + 1


80
    def getFemaleAndMaleData(aggregateUsers):
        aggregateUsersFemale = []
        aggregateUsersMale = []
        if( len(aggregateUsers) > 0 ):
85
            for user in aggregateUsers:
                if( user[1] == 'F' ):
                    aggregateUsersFemale.append(user)
                else:
90                  aggregateUsersMale.append(user)


        return aggregateUsersFemale, aggregateUsersMale


95


    def getHighestMovieRaters(aggregateUsers, count):
```

```
      if( len(aggregateUsers) > 0 and count > 0 ):
100       raterIDRatedMoviesCount = []

          for rater in aggregateUsers:

              #raterTuple: rater id, number of movies rated
105           raterTuple = (rater[0], len(rater[3]))
              raterIDRatedMoviesCount.append(raterTuple)

          raterIDRatedMoviesCount = sorted(raterIDRatedMoviesCount, key=lambda tup:
              tup[1], reverse=True)

110       i = 1
          for rater in raterIDRatedMoviesCount:
              print rater[0], rater[1]

              if( i == count ):
115               break
              i = i + 1



aggregateMovies, aggregateUsers, movies = aggregateMovieAndUserData()
120
aggregateUsersFemale, aggregateUsersMale = getFemaleAndMaleData(aggregateUsers)

getHighestMovieRaters(aggregateUsers, 5)
```

Figure 6: ques6.py at work

# Problem 7

7. Which 5 raters most agreed with each other? Show the raters' IDs and Pearson's r, sorted by r.

Some 5 raters that mostly agree with each other, represented by the IDs are: '100', '105', '107', '111' and '112'

Listing 7: Python script solving problem 7

```python
from math import sqrt
import os, sys




def loadMovieLens(path='/home/vnwala/ml-100k/'):
  # Get movie titles
    movies = {}
    for line in open(path + 'u.item'):
        (id, title) = line.split('|')[0:2]
        movies[id] = title
  # Load data
    prefs = {}
    for line in open(path + 'u.data'):
        (user, movieid, rating, ts) = line.split('\t')
        prefs.setdefault(user, {})
        #prefs[user][movies[movieid]] = float(rating)
        prefs[user][movieid] = float(rating)
    return prefs, movies

prefs, movies = loadMovieLens()



def sim_pearson(prefs,p1,p2):
    # Get the list of mutually rated shared_items
    si={}
    for item in prefs[p1]:
        if item in prefs[p2]: si[item]=1

    # Find the number of elements
    n=len(si)

    # if they have no ratings in common, return 0
    if n==0: return 0

    # Add up all the preferences
    sum1=sum([prefs[p1][it] for it in si])
    sum2=sum([prefs[p2][it] for it in si])

    # Sum up the squares
    sum1Sq=sum([pow(prefs[p1][it],2) for it in si])
    sum2Sq=sum([pow(prefs[p2][it],2) for it in si])

    # Sum up the products
```

      

```python
        pSum=sum([prefs[p1][it]*prefs[p2][it] for it in si])

        # Calculate Pearson score
50      num=pSum-(sum1*sum2/n)
        den=sqrt((sum1Sq-pow(sum1,2)/n)*(sum2Sq-pow(sum2,2)/n))
        if den==0: return 0

        r=num/den
55
        return r




60




65
    def topMatches(prefs,person,n=5,similarity=sim_pearson, reverseSimilarityFlag=True
        ):
        scores=[(similarity(prefs,person,other),other) for other in prefs if other!=
            person]

        # Sort the list so the highest scores appear at the top
70      scores.sort()

        if( reverseSimilarityFlag ):
            scores.reverse()

75      return scores[0:n]




80

    def calculateSimilarItems(prefs, similarityMetric, n=10, reverseSimilarityFlag=
        True, transformMatrixFlag=True):
        '''
        Create a dictionary of items showing which other items they are
85      most similar to.
        '''

        result = {}
        # Invert the preference matrix to be item-centric
90
        if( transformMatrixFlag ):
            #with transform: movie top similarity
            itemPrefs = transformPrefs(prefs)
        else:
95          #without transform: user top similarity
            itemPrefs = prefs
```

```
         #c = 0
         for item in itemPrefs:
100
             scores = topMatches(itemPrefs, item, n=n, similarity=similarityMetric,
                 reverseSimilarityFlag=reverseSimilarityFlag)

             result[item] = scores
         return result
105

     def sim_distance(prefs,person1,person2):
         # Get the list of shared_items
         si = {}
110      for item in prefs[person1]:
             if item in prefs[person2]:
                 si[item]=1

         # if they have no ratins in common, return 0
115      if len(si)==0: return 0

         # Add up the squares of all the differences
         sum_of_squares=sum([pow(prefs[person1][item]-prefs[person2][item],2) for item
             in si])

120      return 1/(1+sqrt(sum_of_squares))




125


     userSimilarityMatrix = calculateSimilarItems(prefs=prefs, similarityMetric=
         sim_distance, n=5, reverseSimilarityFlag=False, transformMatrixFlag=False)
     userSimilarityArrayOfTuples = []
     for userId, userAttr in userSimilarityMatrix.items():
130
         totalSimilarity = 0
         similarUsersArray = []
         for scoreAnduserId in userAttr:
             score = scoreAnduserId[0]
135          userId = scoreAnduserId[1]

             similarUsersArray.append(userId)
             totalSimilarity = totalSimilarity + score

140      userTupleData = (userId, similarUsersArray, totalSimilarity )
         userSimilarityArrayOfTuples.append(userTupleData)


     for userTuple in userSimilarityArrayOfTuples:
145
         userId = userTuple[0]
```

```
        userSimilarItems = userTuple[1]
        totalSim = userTuple[2]


        #any one of these qualifies
        if( totalSim == 0 ):
            print userId, userSimilarItems
```

Figure 7: ques7.py at work



```
112 ['100', '105', '107', '111', '112']
120 ['106', '114', '115', '12', '120']
172 ['124', '138', '148', '156', '172']
138 ['114', '115', '122', '124', '138']
148 ['114', '122', '124', '138', '148']
132 ['114', '118', '122', '124', '132']
212 ['122', '172', '187', '208', '212']
146 ['111', '122', '124', '133', '146']
132 ['114', '115', '122', '124', '132']
212 ['114', '122', '196', '208', '212']
261 ['155', '191', '240', '241', '261']
273 ['105', '111', '147', '171', '273']
124 ['111', '114', '118', '122', '124']
147 ['111', '112', '133', '146', '147']
228 ['122', '172', '208', '212', '228']
124 ['101', '108', '109', '118', '124']
131 ['105', '107', '111', '129', '131']
341 ['166', '238', '273', '34', '341']
162 ['124', '138', '148', '156', '162']
415 ['217', '245', '289', '35', '415']
700 ['124', '366', '471', '571', '700']
778 ['172', '31', '565', '571', '778']
148 ['114', '122', '124', '138', '148']
112 ['100', '105', '107', '111', '112']
191 ['111', '140', '147', '171', '191']
35 ['111', '147', '273', '319', '35']
547 ['289', '341', '36', '519', '547']
19 ['110', '124', '170', '179', '19']
98 ['208', '565', '713', '866', '98']
375 ['187', '208', '281', '366', '375']
208 ['114', '118', '156', '172', '208']
522 ['114', '122', '172', '208', '522']
842 ['260', '317', '335', '418', '842']
```

# Problem 8

8. Which 5 raters most disagreed with each other (negative correlation)? Show the raters' IDs and Pearson's r, sorted by r.

Listing 8: Python script solving problem 8

```python
from math import sqrt
import os, sys



def loadMovieLens(path='/home/vnwala/ml-100k/'):
  # Get movie titles
    movies = {}
    for line in open(path + 'u.item'):
        (id, title) = line.split('|')[0:2]
        movies[id] = title
  # Load data
    prefs = {}
    for line in open(path + 'u.data'):
        (user, movieid, rating, ts) = line.split('\t')
        prefs.setdefault(user, {})
        #prefs[user][movies[movieid]] = float(rating)
        prefs[user][movieid] = float(rating)
    return prefs, movies

prefs, movies = loadMovieLens()



def sim_pearson(prefs,p1,p2):
    # Get the list of mutually rated shared_items
    si={}
    for item in prefs[p1]:
        if item in prefs[p2]: si[item]=1

    # Find the number of elements
    n=len(si)

    # if they have no ratings in common, return 0
    if n==0: return 0

    # Add up all the preferences
    sum1=sum([prefs[p1][it] for it in si])
    sum2=sum([prefs[p2][it] for it in si])

    # Sum up the squares
    sum1Sq=sum([pow(prefs[p1][it],2) for it in si])
    sum2Sq=sum([pow(prefs[p2][it],2) for it in si])

    # Sum up the products
    pSum=sum([prefs[p1][it]*prefs[p2][it] for it in si])
```

```
      # Calculate Pearson score
50    num=pSum-(sum1*sum2/n)
      den=sqrt((sum1Sq-pow(sum1,2)/n)*(sum2Sq-pow(sum2,2)/n))
      if den==0: return 0

      r=num/den
55
      return r




60




65
  def topMatches(prefs,person,n=5,similarity=sim_pearson, reverseSimilarityFlag=True
      ):
      scores=[(similarity(prefs,person,other),other) for other in prefs if other!=
          person]

      # Sort the list so the highest scores appear at the top
70    scores.sort()

      if( reverseSimilarityFlag ):
          scores.reverse()

75    return scores[0:n]




80

  def calculateSimilarItems(prefs, similarityMetric, n=10, reverseSimilarityFlag=
      True, transformMatrixFlag=True):
      '''
      Create a dictionary of items showing which other items they are
85    most similar to.
      '''

      result = {}
      # Invert the preference matrix to be item-centric
90
      if( transformMatrixFlag ):
          #with transform: movie top similarity
          itemPrefs = transformPrefs(prefs)
      else:
95        #without transform: user top similarity
          itemPrefs = prefs
```

```
          #c = 0
          for item in itemPrefs:
100
              scores = topMatches(itemPrefs, item, n=n, similarity=similarityMetric,
                  reverseSimilarityFlag=reverseSimilarityFlag)

              result[item] = scores
          return result
105


      def sim_distance(prefs,person1,person2):
          # Get the list of shared_items
          si = {}
110       for item in prefs[person1]:
              if item in prefs[person2]:
                  si[item]=1

          # if they have no ratins in common, return 0
115       if len(si)==0: return 0

          # Add up the squares of all the differences
          sum_of_squares=sum([pow(prefs[person1][item]-prefs[person2][item],2) for item
              in si])

120       return 1/(1+sqrt(sum_of_squares))




125




      #method 1: distance metric
130   userSimilarityMatrix = calculateSimilarItems(prefs=prefs, similarityMetric=
          sim_distance, n=5, reverseSimilarityFlag=True, transformMatrixFlag=False) #
          reverseSimilarityFlag=True: largest to smallest
      #userSimilarityMatrix = calculateSimilarItems(prefs=prefs, n=5,
          reverseSimilarityFlag=True, transformMatrixFlag=False, similarity=sim_pearson)

      #print len(userSimilarityMatrix)
      userSimilarityArrayOfTuples = []
135   for userId, userAttr in userSimilarityMatrix.items():

          totalSimilarity = 0
          similarUsersArray = []
          for scoreAnduserId in userAttr:
140           score = scoreAnduserId[0]
              userId = scoreAnduserId[1]

              similarUsersArray.append(userId)
              totalSimilarity = totalSimilarity + score
145
```

```
       userTupleData = (userId, similarUsersArray, totalSimilarity )
       userSimilarityArrayOfTuples.append(userTupleData)


150  count = 0
     for userTuple in userSimilarityArrayOfTuples:

         userId = userTuple[0]
         userSimilarItems = userTuple[1]
155      totalSim = userTuple[2]


         #any one of these qualifies
         #pick largest totalSim
160
         print userId, userSimilarItems, totalSim

     print 'count: ', count
```

Figure 8: ques8.py at work

# Problem 9

9. What movie was rated highest on average by men over 40? By men under 40?

For men over 40 we have: 1) Great Day in Harlem, A (1994) 5.0 2) Two or Three Things I Know About Her (1966) 5.0 3) Aparajito (1956) 5.0 4) Strawberry and Chocolate (Fresa y chocolate) (1993) 5.0 5) Little Princess, The (1939) 5.0

For men under 40 we have: 1) Entertaining Angels: The Dorothy Day Story (1996) 5.0 2) Letter From Death Row, A (1998) 5.0 3) Hugo Pool (1997) 5.0 4) Leading Man, The (1996) 5.0 5) Quiet Room, The (1996) 5.0

Listing 9: Python script solving problem 9

```python
def aggregateMovieAndUserData(path='/home/vnwala/ml-100k/'):

    try:

        movies = {}
        aggregateMovieData = []
        for line in open(path + 'u.item'):
            (id, title) = line.split('|')[0:2]
            movies[id] = (title, [], -1)

        users = {}
        #populate user and movie data

        for line in open(path + 'u.data'):
            (user, movieid, rating, ts) = line.split('\t')

            user = user.strip()
            movieid = movieid.strip()
            rating = rating.strip()
            ts = ts.strip()

            users.setdefault(user, {})
            #movie title: movies[movieid][0]
            users[user][movieid] = float(rating)#key is movie title

            '''
            if( user in users ):
                users[user][movies[movieid][0]] = float(rating)
            else:
                users[user] = {}
            '''



            #movies[movieid]: (title, [ArrayOfRatings])
            #movies[movieid][0]: title
            #movies[movieid][1]: array of ratings
            movies[movieid][1].append(float(rating))



        #process movie data
        for movieId, tupleData in movies.items():
            averageRating = sum(tupleData[1]) / float(len(tupleData[1]))
```

```
45                    movietuples = (movieId, tupleData[0], averageRating, len(tupleData[1])
                          )
                   aggregateMovieData.append(movietuples)


           aggregateUserData = []
50         for line in open(path + 'u.user'):
               (userId, age, gender, occupation, zipCode) = line.split('|')

               userTuples = (userId, gender, age, users[userId])
               aggregateUserData.append(userTuples)
55

       except:
           exc_type, exc_obj, exc_tb = sys.exc_info()
           fname = os.path.split(exc_tb.tb_frame.f_code.co_filename)[1]
60         print(fname, exc_tb.tb_lineno, sys.exc_info() )
           return



65     return aggregateMovieData, aggregateUserData, movies


   def getFemaleAndMaleData(aggregateUsers):
       aggregateUsersFemale = []
70     aggregateUsersMale = []
       if( len(aggregateUsers) > 0 ):

           for user in aggregateUsers:
               if( user[1] == 'F' ):
75                 aggregateUsersFemale.append(user)
               else:
                   aggregateUsersMale.append(user)


80     return aggregateUsersFemale, aggregateUsersMale




85 def getHighestRatedMoviesByMenOrWomenUnderAge(aggregateUsersFemaleOrMale, count,
       ageLimit, movies):

       if( count > 0 and count < len(aggregateUsersFemaleOrMale) and ageLimit > 0 and
           len(movies) > 0 ):
           tupleOfMovieRatingDictionary = {}
           movieAverageRatingArrayOfTuples = []
90         for user in aggregateUsersFemaleOrMale:
               if( int(user[2]) < ageLimit ):

                   for movie, rating in user[3].items():
```

```python
95                              if ( movie in tupleOfMovieRatingDictionary ):
                                    tupleOfMovieRatingDictionary[movie].append(rating)
                                else:
                                    tupleOfMovieRatingDictionary[movie] = []
                                    tupleOfMovieRatingDictionary[movie].append(rating)
100
            for movie, ratingsArray in tupleOfMovieRatingDictionary.items():
                averageRating = sum(ratingsArray) / float(len(ratingsArray))

                movieRatingTuple = (movie, averageRating)
105             movieAverageRatingArrayOfTuples.append(movieRatingTuple)


            movieAverageRatingArrayOfTuples = sorted(movieAverageRatingArrayOfTuples,
                key=lambda tup: tup[1], reverse=True)

110         i = 1
            for movieData in movieAverageRatingArrayOfTuples:
                print movies[movieData[0]][0], movieData[1]

                if ( i == count ):
115                 break
                i = i + 1


#input:
    #aggregateUsers: [ (user id, gender, Age, {'movie_title': movie rating}) ]
120




125


def getHighestRatedMoviesByMenOrWomenOverAge(aggregateUsersFemaleOrMale, count,
    ageLimit, movies):

    if ( count > 0 and count < len(aggregateUsersFemaleOrMale) and ageLimit > 0 and
        len(movies) > 0 ):
130         tupleOfMovieRatingDictionary = {}
            movieAverageRatingArrayOfTuples = []

            for user in aggregateUsersFemaleOrMale:

135             if ( int(user[2]) > ageLimit ):
                    for movie, rating in user[3].items():

                        if ( movie in tupleOfMovieRatingDictionary ):
                            tupleOfMovieRatingDictionary[movie].append(rating)
140                     else:
                            tupleOfMovieRatingDictionary[movie] = []
                            tupleOfMovieRatingDictionary[movie].append(rating)
```

```python
145            #average ratings
           for movie, ratingsArray in tupleOfMovieRatingDictionary.items():
               averageRating = sum(ratingsArray) / float(len(ratingsArray))

               movieRatingTuple = (movie, averageRating)
150            movieAverageRatingArrayOfTuples.append(movieRatingTuple)

           #sort
           movieAverageRatingArrayOfTuples = sorted(movieAverageRatingArrayOfTuples,
               key=lambda tup: tup[1], reverse=True)

155        i = 1
           for movieData in movieAverageRatingArrayOfTuples:
               print movies[movieData[0]][0], movieData[1]

               if( i == count ):
160                break
               i = i + 1


aggregateMovies, aggregateUsers, movies = aggregateMovieAndUserData()

165 aggregateUsersFemale, aggregateUsersMale = getFemaleAndMaleData(aggregateUsers)

getHighestRatedMoviesByMenOrWomenOverAge(aggregateUsersMale, 5, 40, movies)
print ''
getHighestRatedMoviesByMenOrWomenUnderAge(aggregateUsersMale, 5, 40, movies)
```
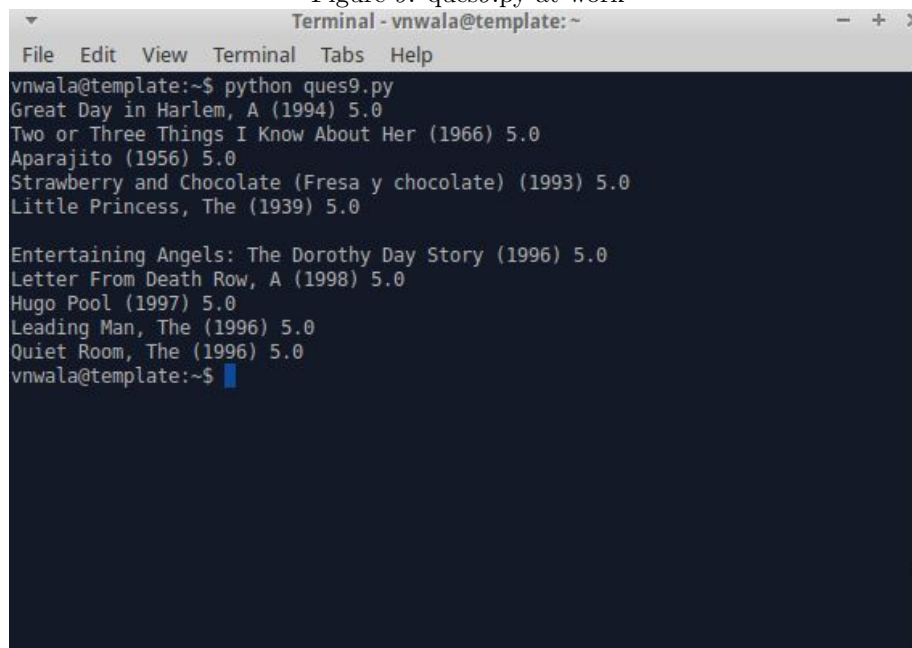
Figure 9: ques9.py at work

# Problem 10

10. What movie was rated highest on average by women over 40? By women under 40?

For women over 40 we have:

1) In the Bleak Midwinter (1995) 5.0 2) Foreign Correspondent (1940) 5.0 3) Swept from the Sea (1997) 5.0 4) Great Dictator, The (1940) 5.0 5) Balto (1995) 5.0

For women under 40 we have: 1) Nico Icon (1995) 5.0 2) Backbeat (1993) 5.0 3) Umbrellas of Cherbourg, The (Parapluies de Cherbourg, Les) (1964) 5.0 4) Everest (1998) 5.0 5) Someone Else's America (1995) 5.0

Listing 10: Python script solving problem 10

```python
def aggregateMovieAndUserData(path='/home/vnwala/ml-100k/'):

    try:

        movies = {}
        aggregateMovieData = []
        for line in open(path + 'u.item'):
            (id, title) = line.split('|')[0:2]
            movies[id] = (title, [], -1)

        users = {}
        #populate user and movie data

        for line in open(path + 'u.data'):
            (user, movieid, rating, ts) = line.split('\t')

            user = user.strip()
            movieid = movieid.strip()
            rating = rating.strip()
            ts = ts.strip()

            users.setdefault(user, {})
            #movie title: movies[movieid][0]
            users[user][movieid] = float(rating)#key is movie title

            '''
            if( user in users ):
                users[user][movies[movieid][0]] = float(rating)
            else:
                users[user] = {}
            '''


            #movies[movieid]: (title, [ArrayOfRatings])
            #movies[movieid][0]: title
            #movies[movieid][1]: array of ratings
            movies[movieid][1].append(float(rating))



        #process movie data
        for movieId, tupleData in movies.items():
            averageRating = sum(tupleData[1]) / float(len(tupleData[1]))
```

      

```
45                    movietuples = (movieId, tupleData[0], averageRating, len(tupleData[1])
                          )
                    aggregateMovieData.append(movietuples)


            aggregateUserData = []
50          for line in open(path + 'u.user'):
                (userId, age, gender, occupation, zipCode) = line.split('|')

                userTuples = (userId, gender, age, users[userId])
                aggregateUserData.append(userTuples)
55

        except:
            exc_type, exc_obj, exc_tb = sys.exc_info()
            fname = os.path.split(exc_tb.tb_frame.f_code.co_filename)[1]
60          print(fname, exc_tb.tb_lineno, sys.exc_info() )
            return



65      return aggregateMovieData, aggregateUserData, movies


    def getFemaleAndMaleData(aggregateUsers):
        aggregateUsersFemale = []
70      aggregateUsersMale = []
        if( len(aggregateUsers) > 0 ):

            for user in aggregateUsers:
                if( user[1] == 'F' ):
75                  aggregateUsersFemale.append(user)
                else:
                    aggregateUsersMale.append(user)


80      return aggregateUsersFemale, aggregateUsersMale




85  def getHighestRatedMoviesByMenOrWomenUnderAge(aggregateUsersFemaleOrMale, count,
        ageLimit, movies):

        if( count > 0 and count < len(aggregateUsersFemaleOrMale) and ageLimit > 0 and
            len(movies) > 0 ):
            tupleOfMovieRatingDictionary = {}
            movieAverageRatingArrayOfTuples = []
90          for user in aggregateUsersFemaleOrMale:
                if( int(user[2]) < ageLimit ):

                    for movie, rating in user[3].items():
```

```python
95                           if( movie in tupleOfMovieRatingDictionary ):
                                 tupleOfMovieRatingDictionary[movie].append(rating)
                             else:
                                 tupleOfMovieRatingDictionary[movie] = []
                                 tupleOfMovieRatingDictionary[movie].append(rating)
100
             for movie, ratingsArray in tupleOfMovieRatingDictionary.items():
                 averageRating = sum(ratingsArray) / float(len(ratingsArray))

                 movieRatingTuple = (movie, averageRating)
105              movieAverageRatingArrayOfTuples.append(movieRatingTuple)


             movieAverageRatingArrayOfTuples = sorted(movieAverageRatingArrayOfTuples,
                 key=lambda tup: tup[1], reverse=True)

110          i = 1
             for movieData in movieAverageRatingArrayOfTuples:
                 print movies[movieData[0]][0], movieData[1]

                 if( i == count ):
115                  break
                 i = i + 1


     #input:
         #aggregateUsers: [ (user id, gender, Age, {'movie_title': movie rating}) ]
120




125


     def getHighestRatedMoviesByMenOrWomenOverAge(aggregateUsersFemaleOrMale, count,
         ageLimit, movies):

         if( count > 0 and count < len(aggregateUsersFemaleOrMale) and ageLimit > 0 and
             len(movies) > 0 ):
130          tupleOfMovieRatingDictionary = {}
             movieAverageRatingArrayOfTuples = []

             for user in aggregateUsersFemaleOrMale:

135              if( int(user[2]) > ageLimit ):
                     for movie, rating in user[3].items():

                         if( movie in tupleOfMovieRatingDictionary ):
                             tupleOfMovieRatingDictionary[movie].append(rating)
140                      else:
                             tupleOfMovieRatingDictionary[movie] = []
                             tupleOfMovieRatingDictionary[movie].append(rating)
```

```
145             #average ratings
            for movie, ratingsArray in tupleOfMovieRatingDictionary.items():
                averageRating = sum(ratingsArray) / float(len(ratingsArray))

                movieRatingTuple = (movie, averageRating)
150             movieAverageRatingArrayOfTuples.append(movieRatingTuple)

            #sort
            movieAverageRatingArrayOfTuples = sorted(movieAverageRatingArrayOfTuples,
                key=lambda tup: tup[1], reverse=True)

155         i = 1
            for movieData in movieAverageRatingArrayOfTuples:
                print movies[movieData[0]][0], movieData[1]

                if( i == count ):
160                 break
                i = i + 1


    aggregateMovies, aggregateUsers, movies = aggregateMovieAndUserData()

165 aggregateUsersFemale, aggregateUsersMale = getFemaleAndMaleData(aggregateUsers)


    getHighestRatedMoviesByMenOrWomenOverAge(aggregateUsersFemale, 5, 40, movies)
    print ''
    getHighestRatedMoviesByMenOrWomenUnderAge(aggregateUsersFemale, 5, 40, movies)
```
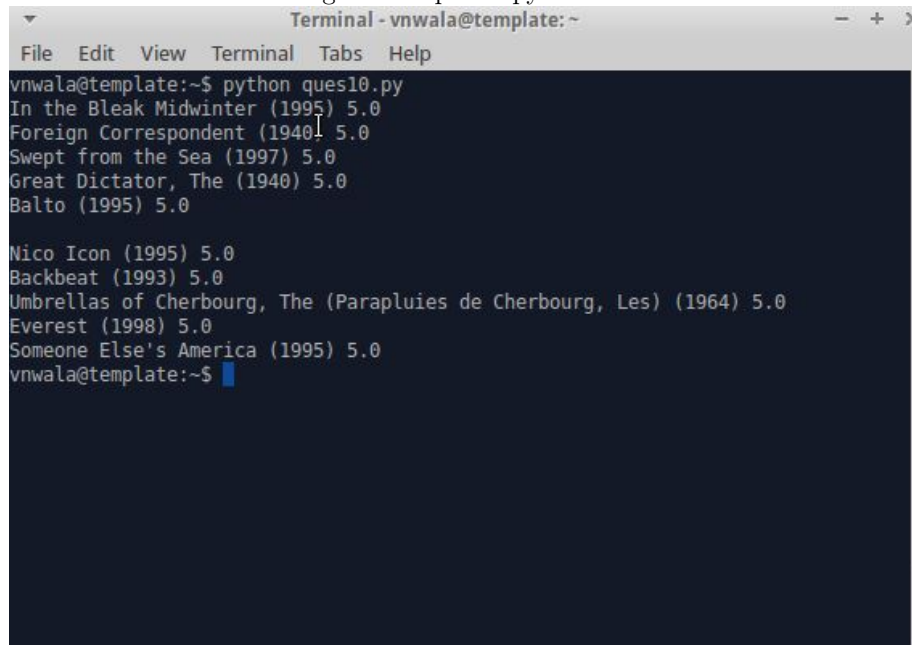
Figure 10: ques10.py at work

# Conclusion

To conclude, I should state that Alexander Nwala was a huge contributor for my answers from question 6 to 10, so those answers will in some or most cases have the same syntatic and functional properties similar to his code. Some questions were answered in part, depending on what I was able to do.

# References

[1] arthur e. Programming-collective-intelligence. https://github.com/arthur-e/Programming-Collective-Intelligence/blob/master/chapter2/recommendations.py, 24 Decenber 2012.

[]